

Module 1: Nonparametric Preliminaries

Selecting Smoothing Parameters

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 8th, 2014

©Emily Fox 2014

1

Smoothing Parameter

- In both ridge and lasso regression, we saw that the parameter λ controlled the solution
 - Often, can straightforwardly equate with effective degrees of freedom

- Which λ (\rightarrow estimator) should we choose???

Want good predictions

↑
Linear smoothers
 $\gamma_\lambda = \text{tr}(L_\lambda)$
↑
"hat matrix"

©Emily Fox 2014

2

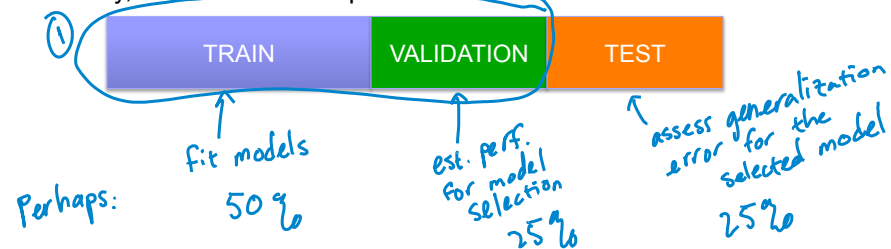
Two Goals

- ① **Model Selection:** estimating the performance of models in order to select the best one

- E.g., choosing λ

- ② **Model Assessment:** having chosen a final model, estimate its prediction error (generalization error) on new data

- Ideally, divide data into 3 parts



©Emily Fox 2014

3

Focus on Model Selection

- Which estimator/smoothing parameter should we choose?



- Recall metrics for assessing the performance of an estimator...

©Emily Fox 2014

4

Measuring Predictive Performance

- Assume estimate $\hat{f}_n(\cdot)$ based on training data y_1, \dots, y_n
fixed
- The **generalization error** provides a measure of predictive performance

$$GE(\hat{f}_n) = E_{Y, X} [L(Y, \hat{f}_n(X))]$$
fixed

©Emily Fox 2014

5

Measuring Predictive Performance

- Assume L_2 loss $Y = f(x) + \epsilon$ $E[\epsilon] = 0$ $\text{var}(\epsilon) = \sigma^2$
- Averaging over repeat training sets $\mathbf{Y}_n = Y_1, \dots, Y_n$ we get the **predictive risk** at x^*

$$E_{Y^*, \mathbf{Y}_n} [(Y^* - \hat{f}_n(x^*))^2] = E_{Y^*, \mathbf{Y}_n} [(Y^* - f(x^*) + f(x^*) - \hat{f}_n(x^*))^2]$$

test training scn of training data

$$= E_{Y^*} [(Y^* - f(x^*))^2] + E_{\mathbf{Y}_n} [(\hat{f}_n(x^*) - f(x^*))^2] + 2 E_{Y^*, \mathbf{Y}_n} [(Y^* - f(x^*))(\hat{f}_n(x^*) - f(x^*))]$$

0

$$= \sigma^2 + \text{MSE}(\hat{f}_n(x^*))$$

"irreducible error" "risk"

- Recall $\text{MSE}[\hat{f}_n(x)] = \text{bias}(\hat{f}_n(x))^2 + \text{var}(\hat{f}_n(x))$

©Emily Fox 2014

6

Measuring Predictive Performance

- Finally, let's average over covariates x *(focus on MSE bc can't avoid σ^2)*

- Integrated MSE** $\int \text{MSE}(\hat{f}_n(x)) p(x) dx$
summary over all inputs

- Average MSE** $\frac{1}{n} \sum_{i=1}^n \text{MSE}(\hat{f}_n(x_i))$
Monte Carlo est: $x_i \sim P \quad i=1, \dots, n$

- Note: **avg. pred. risk = $\sigma^2 + \text{avg. MSE}$** *still our focus*

$$\frac{1}{n} \sum_{i=1}^n E_{Y_i^*, Y_n} [(Y_i^* - \hat{f}_n(x_i))^2]$$
training data
new obs. at x_i

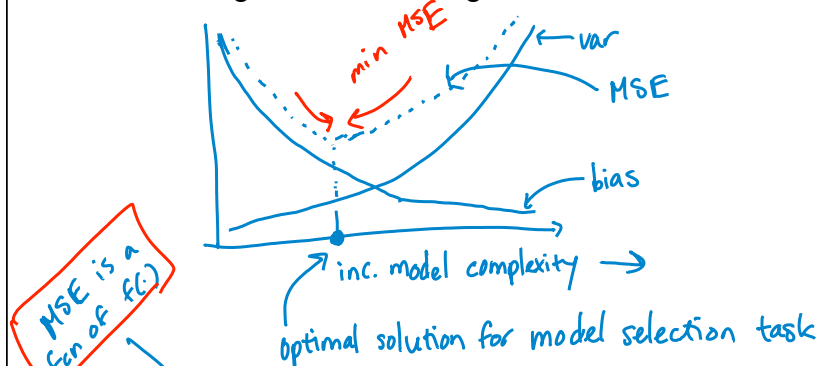
©Emily Fox 2014

7

Bias-Variance Tradeoff

recall polynomial reg. example

- Minimizing risk = balancing bias and variance



- Note: $f(x)$ is unknown, so cannot actually compute MSE

©Emily Fox 2014

8

Focus on Model Selection

- Which estimator/smoothing parameter should we choose?



- We saw that minimizing (average) prediction error can be equated with minimizing (average) MSE
- With a validation set, we can estimate the prediction error

$$\frac{1}{m} \sum_{i=1}^m (y_i^* - \hat{f}_n^\lambda(x_i^*))^2$$

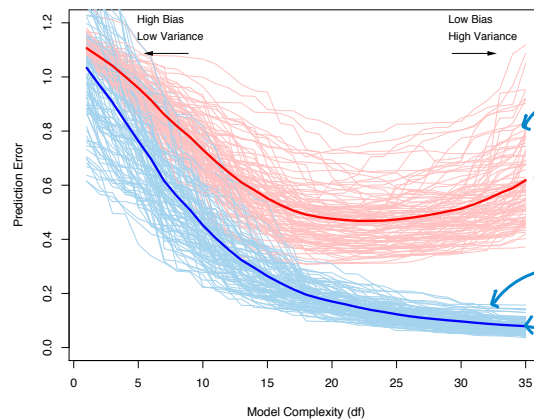
Handwritten annotations: An arrow points from the text 'est. from training data of size n using λ ' to the \hat{f}_n^λ term. Another arrow points from the text 'obs. in validation set of size m' to the y_i^* term.

©Emily Fox 2014

9

In Practice...

- Minimizing risk = balancing bias and variance



From Hastie, Tibshirani, Friedman

©Emily Fox 2014

10

Data Scarce Approximations

- Often, we do not have enough data to form suitably sized training and validation sets
 - What is a good training/test split? Sensitivity?
 - Typically want to use as much data for training as possible
- Rely on other approximations *using in-sample data*

©Emily Fox 2014

11

~~DO NOT USE THIS~~ Approx 1: Training Data Only

- **Goal:** Minimize average MSE

$$\min_{\lambda} E \left[\frac{1}{n} \sum_{i=1}^n (f(x_i) - \hat{f}_n^{\lambda}(x_i))^2 \right]$$

- **Solution:** Use training error

$$\min \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_n^{\lambda}(x_i))^2$$

↑ training obs.

$$\text{training error} = \frac{RSS}{n}$$

BAD

biased downwards + leads to overfitting (undersmoothing)

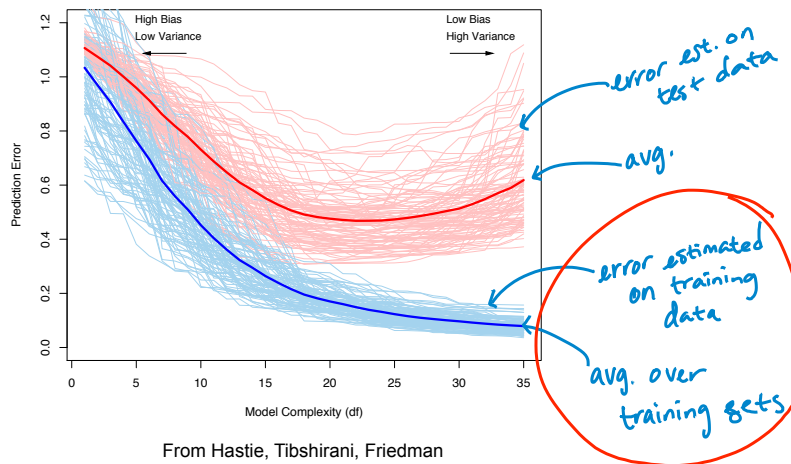
Data was used twice → est. fcn ← tried to min $\frac{RSS}{n}$,
 → est. risk so under est. risk

©Emily Fox 2014

12

In Practice...

- Minimizing risk = balancing bias and variance



©Emily Fox 2014

13

Approx 2: Cross Validation

- **Goal:** Minimize average MSE

$$\min_{\lambda} E \left[\frac{1}{n} \sum_{i=1}^n (f(x_i) - \hat{f}_n^{\lambda}(x_i))^2 \right]$$

- **Solution:** Mimic heldout data using *training* data

- Leave-one-out (LOO) cross validation (CV) algorithm:

- Estimate fit using all but i^{th} data point $\hat{f}_{-i}^{\lambda} \leftarrow \text{no obs. } y_i$
- Predict i^{th} observation
- Repeat for all i

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_{-i}^{\lambda}(x_i))^2$$

- Repeat for all values of λ

©Emily Fox 2014

14

Approx 2: Cross Validation

- Reasoning

$$CV = E[(Y_i - \hat{f}_{-i}^\lambda(x_i))^2] = E[(Y_i - f(x_i) + f(x_i) - \hat{f}_{-i}^\lambda(x_i))^2]$$

$$= \sigma^2 + E[(f(x_i) - \hat{f}_{-i}^\lambda(x_i))^2]$$

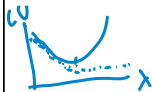
$$\approx \sigma^2 + E[(f(x_i) - \hat{f}_n^\lambda(x_i))^2]$$

- For linear smoothers

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - L_{ii}^\lambda} \right)^2$$

only do fit once (per λ)

\nwarrow i th diag. element of hat matrix using λ



- Warning: Curves can be very flat...Don't just choose and use without thinking. Some rules of thumb (see Elements of Statistical Learning)

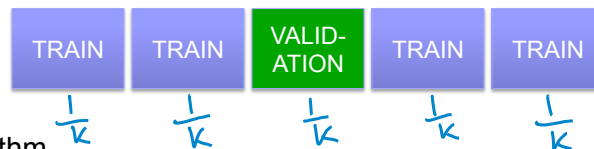
©Emily Fox 2014

15

Approx 2: Cross Validation

- K-fold cross validation

typically $k=5, 10$



- Algorithm

- Fit model using data with k th fraction removed
- Using fitted model, compute

$$CV_k(\lambda) = \frac{1}{n_k} \sum_{i \in J(k)} (y_i - \hat{f}_{-k}^\lambda(x_i))$$

\nwarrow w/o k th block

\nwarrow indices for k th block

- Store

$$CV(\lambda) = \frac{1}{K} \sum_{k=1}^K CV_k(\lambda)$$

- Repeat for each value of λ using same split of the data

©Emily Fox 2014

16

Approx 3: Generalized CV

- Recall LOO ordinary CV for linear smoothers

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - L_{ii}^\lambda} \right)^2$$

- Instead of L_{ii}^λ , use $\frac{1}{n} \sum_{i=1}^n L_{ii}^\lambda = \frac{1}{n} \text{tr}(L^\lambda) = \frac{\nu_\lambda}{n}$

$$GCV(\lambda) = \frac{1}{n} \sum \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - \frac{\nu_\lambda}{n}} \right)^2$$

- Often very close to OCV solution

©Emily Fox 2014

17

Approx 3: Generalized CV

$$GCV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - \frac{\nu_\lambda}{n}} \right)^2$$

- One motivation: Invariance to orthonormal transformations

$$y, x \rightarrow Qy, Qx \quad \text{where} \quad Q^T Q = Q Q^T = I$$

$$\min (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

$$\min (Qy - QX\beta)^T (Qy - QX\beta) + \lambda \beta^T \beta$$

If L^λ is a linear smoother for original data, $L_Q^\lambda = Q L^\lambda Q^T$ is for trans. data

$$\text{tr}(L_Q^\lambda) = \text{tr}(Q L^\lambda Q^T) = \text{tr}(L^\lambda Q^T Q) = \text{tr}(L^\lambda)$$

GCV scores will be the same

©Emily Fox 2014

18

Approx 3: Generalized CV

$$GCV(\lambda) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}_n^\lambda(x_i)}{1 - \frac{\nu_\lambda}{n}} \right)^2$$

$(1-x)^{-2}$

- Using $(1-x)^{-2} \approx 1 + 2x$

$$GCV(\lambda) \approx \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_n^\lambda(x_i))^2 + 2 \frac{\nu_\lambda}{n} \left(\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_n^\lambda(x_i))^2 \right)$$

$\hat{\sigma}^2$

\approx Mallows's C_p stat

(but not exactly the right $\hat{\sigma}^2$)

©Emily Fox 2014

19

Approx 4: Mallows C_p Statistic

- Goal:** Minimize average MSE

$$\min_{\lambda} E \left[\frac{1}{n} \sum_{i=1}^n (f(x_i) - \hat{f}_n^\lambda(x_i))^2 \right]$$

for linear smoothers

- Solution:** Approximate directly

$$\begin{aligned} \text{avg. MSE} &= \frac{1}{n} E \left[(f - \hat{f}_n^\lambda)^T (f - \hat{f}_n^\lambda) \right] = \frac{1}{n} E \left[(Y - \epsilon - L^\lambda Y)^T (Y - \epsilon - L^\lambda Y) \right] \\ &= \frac{1}{n} E \left[(Y - L^\lambda Y)^T (Y - L^\lambda Y) \right] - \sigma^2 + \frac{2}{n} \nu_\lambda \sigma^2 \end{aligned}$$

\hat{f}_n^λ

$$\begin{aligned} \text{uses } E[\epsilon^T L^\lambda \epsilon] &= E[\text{tr}(\epsilon^T L^\lambda \epsilon)] = E[\text{tr}(L^\lambda \epsilon \epsilon^T)] \\ &= \text{tr}(L^\lambda I \sigma^2) = \sigma^2 \nu_\lambda \end{aligned}$$

©Emily Fox 2014

20

Approx 4: Mallows C_p Statistic

$$\text{avg. MSE} = \frac{1}{n} E[(Y - L^\lambda Y)^T (Y - L^\lambda Y)] - \sigma^2 + \frac{2}{n} \nu_\lambda \sigma^2$$

- Estimate avg. MSE as

$$\frac{1}{n} \text{RSS}^\lambda - \frac{1}{n} (n - 2\nu_\lambda) \hat{\sigma}_{\max}^2$$

\Updownarrow
 min Mallows's C_p
 $\frac{\text{RSS}^\lambda}{\hat{\sigma}_{\max}^2} - (n - 2\nu_\lambda)$

using a maximal model

- Note: Arises from considering L_2 loss. Log-likelihood loss leads to AIC. For BIC, consider Bayesian model selection

©Emily Fox 2014

21

Bayesian Model Selection

- Assume some M possible models

- Model M_m $m=1, \dots, M$ has parameters θ_m and prior $p(\theta_m | M_m)$
- Prior over models $p(M_m)$

- Model posterior

$$p(M_m | Z) \propto p(M_m) p(Z | M_m)$$

$$\propto p(M_m) \int p(Z | \theta_m, M_m) p(\theta_m | M_m) d\theta_m$$

- Compare models:

$$\frac{p(M_m | Z)}{p(M_\ell | Z)} = \frac{p(M_m) p(Z | M_m)}{p(M_\ell) p(Z | M_\ell)} \geq 1$$

Posterior odds

Often, uniform prior

Bayes factor

eg mean + cov of Gauss or β_j

©Emily Fox 2014

22

Bayesian Model Selection

- For Bayes factor, approximate

$$\log p(Z | M_m) \approx \log p(Z | \hat{\theta}_m, M_m) - \frac{\nu_m}{2} \log n + O(1)$$

Handwritten notes: "Laplace + ..." above the equation, "ML est." with an arrow pointing to $\hat{\theta}_m$, and "# of free params" with an arrow pointing to ν_m .

- If loss is $-2 \log p(Z | \hat{\theta}_m, M_m)$, then equivalent to BIC
 - Minimizing BIC = maximizing approximated posterior

- However, in addition to being able to select the best model, in Bayesian framework we also get the relative merit of each

$$\approx \frac{e^{-\frac{1}{2} \text{BIC}_m}}{\sum_{\ell=1}^M e^{-\frac{1}{2} \text{BIC}_\ell}}$$

- BIC is asymptotically consistent, but AIC is not
- For finite samples, BIC tends to choose too simple models

©Emily Fox 2014

23

Reading

- Hastie, Tibshirani, Friedman: 7.2 (again), 7.4-7.7, 7.10
- Wakefield: 10.6 (up to 10.6.4)

©Emily Fox 2014

24

What you should know...

- Model selection vs. model assessment tasks
- Training/validation/test split
- In-sample approaches for selecting the smoothing parameters:
 - Training error = BAD
 - Cross validation (CV)
 - LOO
 - K-fold
 - Generalized cross validation (GCV)
 - Mallows's C_p
- Bayesian model selection

©Emily Fox 2014

25

Module 2: Splines and Kernel Methods

Spline Model Overview,
Regression Splines,
Smoothing Splines

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 8th, 2014

©Emily Fox 2014

26

Moving Beyond Linearity

- So far we have assumed standard linear models

$$\min_{\beta} \|y - X\beta\|_2^2 \quad \leftarrow f(x) = \beta^T x$$

- In the case of many predictors relative to number of observations, we considered penalized regression to avoid overfitting

$$\min_{\beta} \|y - X\beta\|_2^2 + \lambda \|\beta\|$$

- Often a convenient form, and necessary to assume simple structure to avoid overfitting in data-scarce regimes, but linear assumption rarely holds in practice

©Emily Fox 2014

27

Moving Beyond Linearity

- Consider generic functional forms (univariate x for now)

$$\min_f \|y - f(x)\|_2^2$$

- If constrained to linear forms \rightarrow LS soln
- If arbitrary \rightarrow interpolator ... overfitting

- As before, penalize complexity. Here, in terms of roughness.

$$\min_f \|y - f(x)\|_2^2 + \lambda \int f''(x)^2 dx$$

- If $\lambda \rightarrow 0$, interpolator
- If $\lambda \rightarrow \infty$, LS soln (line ... no 2nd der.)

- Remarkable result: Explicit, finite-dimensional minimizer

TBD natural cubic spline w/ knots at data pts
"smoothing spline"

©Emily Fox 2014

28

Backtrack a bit...

- Instead of just considering input variables x (potentially mult.), augment/replace with transformations = “input features”

- **Linear basis expansions** maintain linear form in terms of these transformations

$$f(x) = \sum_{m=1}^M \beta_m h_m(x)$$

← trans. ← linear in these transformations

- What transformations should we use?

- $h_m(x) = x_m \rightarrow$ linear model
- $h_m(x) = x_j^2, \quad h_m(x) = x_j x_k \rightarrow$ polynomial reg.
- $h_m(x) = I(L_m \leq x_k \leq U_m) \rightarrow$ piecewise constant
- ...

©Emily Fox 2014

29

Piecewise Polynomial Fits

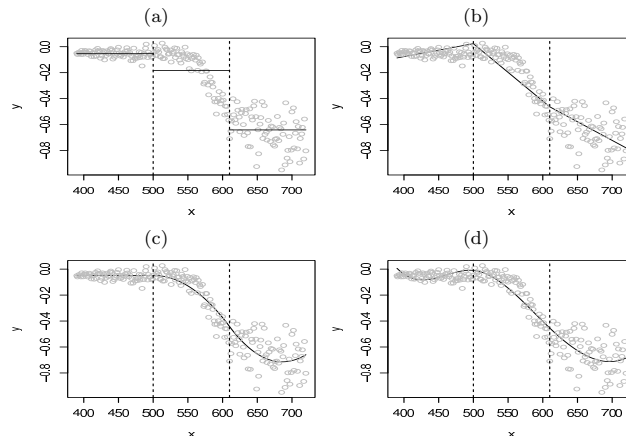
- Again, assume x univariate
- Polynomial fits are often good locally, but not globally
 - Adjusting coefficients to fit one region can make the function go wild in other regions
- Consider **piecewise polynomial** fits
 - Local behavior can often be well approximated by low-order polynomials

©Emily Fox 2014

30

Piecewise Polynomial Fits

LIDAR Data Example



From Wakefield book

©Emily Fox 2014

31

Piecewise Constant/Linear Fits

■ Example 1: Piecewise constant, with 3 basis functions

$$h_1(x) = \mathbb{I}(x \leq \xi_1) \quad \text{"knot"}$$

$$h_2(x) = \mathbb{I}(\xi_1 \leq x \leq \xi_2)$$

$$h_3(x) = \mathbb{I}(\xi_2 \leq x)$$

■ Resulting model:
$$f(x) = \sum_{m=1}^3 \beta_m h_m(x)$$

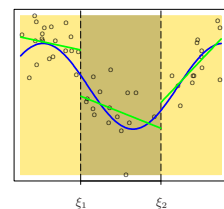
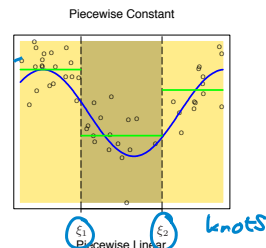
- Fit: Take mean of data in each region

$$\hat{\beta}_m = \bar{y}_m \leftarrow \text{data in region}$$

■ Example 2: Piecewise linear

- Add three basis functions:

$$h_{m+3} = h_m(x)x \quad m=1,2,3$$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

32

Regression Splines – Linear

- Resulting piecewise linear model:

$$f(x) = I(x < \xi_1)(\beta_1 + \beta_4 x) + I(\xi_1 \leq x < \xi_2)(\beta_2 + \beta_5 x) + I(\xi_2 \leq x)(\beta_3 + \beta_6 x)$$

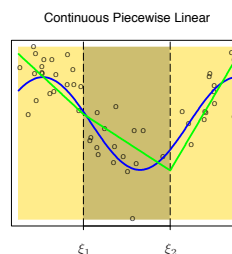
of params? 6

- Typically prefer continuity...

Enforce $f(\xi_1^-) = f(\xi_1^+)$
 $f(\xi_2^-) = f(\xi_2^+)$

Which implies
 $\beta_1 + \beta_4 \xi_1 = \beta_2 + \beta_5 \xi_1$
 $\beta_2 + \beta_5 \xi_2 = \beta_3 + \beta_6 \xi_2$

params?
 $6 - 2 = 4$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

33

Regression Splines – Linear

- More directly, we can use the **truncated power basis**

$$h_1(x) = 1$$

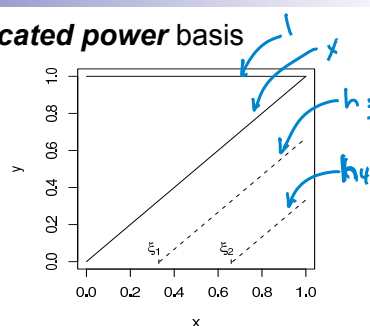
$$h_2(x) = x$$

$$h_3(x) = (x - \xi_1)_+$$

$$h_4(x) = (x - \xi_2)_+$$

- Resulting model:

$$f(x) = \beta_0 + \beta_1 x + \beta_2 (x - \xi_1)_+ + \beta_3 (x - \xi_2)_+$$



From Wakefield book

- Continuous at the knots because all prior basis functions are contributing to the fit up to any single x

©Emily Fox 2014

34

Reading

- Hastie, Tibshirani, Friedman: 5.1-5.5 (skipping 5.3)
- Wakefield: 11.1.1-11.2.3

What you should know...

- Linear basis expansions
- Regression splines
 - Cubic splines, natural cubic splines, ...
 - Interpretation as a linear smoother
 - Degrees of freedom
- Smoothing splines
 - Arising from penalized regression setting with smoothness penalty
 - Cubic spline basis with knots at every data point